

Chunking-Coordinated-Synthetic Approaches to Large-Scale Kernel Machines

Francisco J. González-Castaño¹ and Robert R. Meyer²

¹ Departamento de Ingeniería Telemática, Universidad de Vigo, Spain
ETSI Telecomunicación, Campus, 36200 Vigo, Spain

² Computer Sciences Department, University of Wisconsin-Madison, USA
javier@det.uvigo.es, rrm@cs.wisc.edu

Abstract. We consider a kernel-based approach to nonlinear classification that coordinates the generation of “synthetic” points (to be used in the kernel) with “chunking” (working with subsets of the data) in order to significantly reduce the size of the optimization problems required to construct classifiers for massive datasets. Rather than solving a single massive classification problem involving all points in the training set, we employ a series of problems that gradually increase in size and which consider kernels based on small numbers of synthetic points. These synthetic points are generated by solving and combining the results of relatively small nonlinear unconstrained optimization problems. In addition to greatly reducing optimization problem size, the procedure that we describe also has the advantage of being easily parallelized. Computational results show that our method efficiently generates high-performance simple classifiers on a problem involving a realistic dataset.

1 Introduction

Suppose that the following classification problem is given:

A set of m training points in a n -dimensional space is given by the $m \times n$ matrix Y . These points are divided into two classes, types 1 and -1. A classifier is to be constructed using this data and a nonlinear kernel function K that is a mapping from $R^n \times R^n$ to R . We assume that construction of the corresponding classifier involves generating a function $g(x) = wK(C, x)$, where w is a row vector of size s and C is a set of s points in R^n and it is understood that, for such a set of points C , $K(C, x)$ is a column vector of size s whose entries are given by $K(c_i, x)$ for the rows c_i of C (in a further extension of this notation below, we will assume that for any two matrices A and B each of which has n columns, that $K(A, B)$ is a matrix whose (i, j) entry is $K(A_i, B_j)$ where A_i and B_j are the rows corresponding to indices i and j . In addition, we will use a juxtaposition of vectors and matrices to indicate a product operation, without introducing an extra symbol to indicate the transposition of the vector that may be required). The points used in a set C will be termed classifier points. While it is customary to set $C = Y$, one of the goals of this paper is to investigate alternative approaches to constructing C , particularly when Y is a large dataset.

In conjunction with the determination of the weights w described above, a constant γ is also computed so that the resulting classifier designates (as correctly as possible) a point as type 1 if $g(x) > \gamma$ and a point as type -1 if $g(x) < \gamma$.

We assume the kernel *Mercer inner product property* [2] $K(y, x) = f(y)f(x)$, where f is a mapping from R^n to another space R^k , and $f(y)f(x)$ represents an inner product (for a simple example, suppose that $n = 2$ and $K(y, x) = (yx)^2$. Then by taking $f(z) = (z_1^2, 2^{1/2}z_1z_2, z_2^2)$, where z is a 2-vector, it is easy to verify that $K(y, x) = f(y)f(x)$ holds for any pair of vectors). The original space R^n is finite dimensional and is referred to as “input space” whereas R^k may have a much higher dimension (or even be infinite-dimensional) and is referred to as “feature space” (in the simple example given, $n = 2$ and $k = 3$, but note that even for this simple quadratic kernel the dimension of the feature space is proportional to the square of the dimension of the input space). Many commonly used nonlinear classifiers (such as Gaussians or homogeneous polynomials) have this Mercer property, and although evaluation of the corresponding function f itself is in general not computationally practical, the inner product representation $K(y, x) = f(y)f(x)$ allows classifier approximation problems to be formulated in a manner that is computationally tractable. In fact, by taking advantage of this Mercer property, we will be able to solve via a relatively small unconstrained nonlinear program (NLP), the problem of approximating a classifier with $C = Y$ (or subsets of Y) by classifiers based on very small numbers of “synthetic points” that are not necessarily in Y .

Classifiers are usually constructed by taking the set C to be the set Y of training points and then solving an optimization problem for the values of w and γ . However, for massive datasets this choice of C results in intractably large optimization problems. Moreover, it may be the case that even for medium-sized datasets, the classifiers using $C = Y$ can be outperformed by alternative choices of C . For example, if the type 1 points form a cluster about a center point z that is not in Y , and the type -1 points lie beyond this cluster, then an ideal classifier may be constructed using a Gaussian kernel involving the point z alone rather than any of the many points from Y . As observed in a noisy dataset below, the generalizability associated with alternative smaller choices of C may also be better than the choice $C = Y$, since the latter can lead to overtraining.

In the initialization step of our method we construct classifiers using C 's that are small subsets of Y , and then in successor steps we approximate classifiers for successively larger subsets of the training set by using C 's corresponding to even smaller sets of synthetic points (as opposed to using points from Y). We describe this procedure in the next section. It should be noted that this chunking strategy of considering successively larger subsets is one of the elements that differentiates this research from that of [2], which emphasizes synthetic point generation as an *a posteriori* procedure applied to simplify a classifier obtained in the standard manner from the full dataset. A second difference is our coordination of the results of parallel optimization problems for the synthetics.

The idea of simplifying classifiers by using synthetic points was introduced in [2]. Alternative approaches for generating synthetic points are described in [3,8].

Other approaches that seek to reduce kernel-classifier complexity are discussed in [7]. In the linear case, classifiers have also been simplified by means of feature selection via concave programming [1].

2 Chunking-Coordinated-Synthetic (CCS) Approaches

Our approaches are based on a sequence of classifier generating problems interspersed with a sequence of classifier approximation problems (that generate synthetic points). The former we refer to as classifier problems and the latter, as approximation problems. We now define the format of the problems used in the specific implementation presented below.

LP(S, C) is the linear programming classifier generating problem:

$$\min_{u, \gamma, E} \quad \nu \|E\|_1 + \|u\|_1 \quad \text{s.t.} \quad D(uK(C, S) - \gamma \mathbf{1}) + E \geq \mathbf{1}, \quad E \geq \mathbf{0}, \quad (1)$$

where ν is a weighting factor, $uK(C, S)$ is the vector obtained by applying the kernel corresponding to $uK(C, \cdot)$ to each element of the subset S of the training set, D is a diagonal matrix of ± 1 's corresponding to the appropriate classification of the corresponding element of S , E is a vector of "errors" (with respect to the "soft" margin bounding surfaces corresponding to $uK(C, x) = \gamma \pm 1$), $\mathbf{1}$ is a vector of 1's, and $\mathbf{0}$ is a vector of 0's. Note that the standard classification problem in this framework would be obtained by setting $C = S = Y$, but this leads to a problem with $O(m)$ constraints and variables, which is intractable if m is large. Other classifier problems (such as classifiers with quadratic objective terms) may be substituted for the LP in (1). In particular, in future research we will experiment with further reductions in problem size via the use of *unconstrained* classifier models.

Using the notation $f(C)$ to denote the array obtained by applying f to each of the s points of C , consider the classifier term $u^*f(C)$, where u^* is obtained by solving (1), and generate an approximation (using $t < s$ points) to this term by considering the problem NLP(u^*, C, t), which is the unconstrained problem $\min_{w, Z} d(u^*f(C), wf(Z))$, where $d()$ provides a measure of distance. As a particular case, in this research we used the 2-norm:

$$\min_{w, Z} \quad \|u^*f(C) - wf(Z)\|_2^2, \quad (2)$$

where w is a vector of size t , and Z is a set of t synthetic points, which may be distinct from points in C . Thus, the optimal solution $w^*f(Z^*)$ approximates the classifier term $u^*f(C)$. In order to avoid computation with f in feature space, this NLP is re-formulated by expanding the squared 2-norm and applying the Mercer inner product property to obtain an equivalent problem expressed in input space [2]. The expanded problem (2) is $\min_{w, Z} u^*f(C)f(C)u^* - 2u^*f(C)f(Z)w + wf(Z)f(Z)w$, so the corresponding problem in input space via the inner product property is $\min_{w, Z} u^*K(C, C)u^* - 2u^*K(C, Z)w + wK(Z, Z)w$. Note that a numerical value of $K(y, x)$ for a specific pair (y, x) is computed by operations in input space (for example, we may evaluate $K(y, x) = (yx)^p$ where p is a positive

integer and yx represents the inner product in input space. Evaluation of $K(y, x)$ via the expression $f(y)f(x)$ would be impractical for $p > 1$ and large n).

The CCS strategy considered here allows different algorithmic formulations. We now describe the i -th stage of the multistage CCS algorithm that we used in this research. In our approach we perform p independent runs of the i -th stage process, using a different random subset of training points in each run; the results of these p runs are then combined as described below to provide the initial set of classifier points C_{i+1} for stage $i + 1$. In our results we use the settings $p = 10$, $t = 10$ (number of points used in classifier approximation problems).

At each successive stage, the size of the subset of the training set used to provide data for the classification problem is increased as described below until a classifier problem for the full training set is reached.

Stage i process (performed p times using p independent samples from the training set):

- a) Let C_i be a set of s points in R^n , with $s = pt \ll m$ (C_0 is chosen as a subset of the training set; for $i > 0$ the C_i are sets of synthetic points obtained from the preceding stage).
- b) Let Y_i be a randomly chosen subset of Y such that $size(Y_i) > size(Y_{i-1})$; solve the **classifier problem** $LP(Y_i, C_i)$.
- c) Letting Z_i, w_i be a set of initial values for a nonlinear optimization procedure, solve the **approximation problem** $NLP(u_i^*, C_i, t)$ to obtain a set of synthetics Z_i^* .
- d) Solve the **classifier problem**: $LP(Y_i, Z_i^*)$.

Coordination Step (coordinates the p i -th stage processes): for each run performed in stage $i + 1$, C_{i+1} is the **union** of all Z_i^* for the p runs of the i -th stage (yielding a total of $s = pt$ potential classifier points for each of the initial classifier problems of step (b) in stage $i + 1$).

It should be noted that step (d) of this chunking strategy serves to validate the choices made for the strategic parameters p and t in the sense that the testing set correctness should be similar for steps (b) and (d) and correctness should stabilize as the subsets of Y approach the size of the original training set. This behavior was observed in the computational results that we now present.

Important remark: alternative implementations of step (c) could employ other strategies to generate appropriate synthetic points. Some possible methods are described in [3,8,10].

3 Computational Results, Modified USPS Problem

In order to evaluate the CCS algorithm, we considered a variant of the US Postal Service problem [9]. The USPS problem is composed of 9,298 patterns, 256 features each, corresponding to scanned handwritten numbers (0-9). Thus, it is a multi-category problem. Roughly, each category has the same number of representatives. To generate a (harder) modified version, inspired by [8], we add noise to a sample of USPS points in order to simulate a noisier dataset

with less separability. We use a USPS sample as centers. In order to balance the problem, the sample is composed of 90 random '8' patterns and a random 10-pattern subset for each of the other nine digits. The resulting 180-center set was used as input for the NDC data generator [6] (thus, NDC centers are real data, instead of random points). The NDC generator typically assigns labels to centers by means of a random separating hyperplane, but in our case, this assignment method is not necessary, because labels are known *a priori*: 1 for '8' centers and -1 for the other centers and (as is usual in NDC) all points that are generated from a specific center are assigned the center's label. The overall result is a training set with approximately 4000 points of category 1 and 4000 of category -1. A Gaussian kernel $K(y, x) = \exp(-|y - x|^2 / 128)$ was used as in [8], where it provided good results. All our runs took place on a node of a Sun Enterprise machine with an UltraSPARC II processor at 250 MHz. Step (b) and any other linear problems in our tests were solved using CPLEX 6.0. The nonlinear approximation problem in step (c) was solved by invoking CONOPT 2 via GAMS 2.50A (we chose an "off-the-shelf" NLP solver for this initial implementation, and found it to be quite effective. Other NLP techniques could, of course, be more efficient for this particular class of applications. With CONOPT 2 we found that the solution process was accelerated by imposing box constraints on the variables w . In particular, we impose the constraints: $-\alpha|u_i^*| \leq w \leq \alpha|u_i^*|$ for some $\alpha > 0$. The CONOPT 2 option `rtredg`, reduced gradient tolerance for optimality, was set to 1e-4). Both CPLEX and GAMS were called from a C program that implemented the whole algorithm.

All settings were NDC defaults, except for the number of points (8000 used here), number of features (256) and expansion factor (25) (this factor controls the expansion of point clouds around centers, using a random covariance matrix for each center). The expansion factor was selected to produce a separability in the range of 70% for the Gaussian kernel $K(y, x) = \exp(-|y - x|^2 / 128)$ in a standard 1-norm classification problem. The value $\nu = 10.0$ was found to yield the best training correctness for this dataset. Four stages were used, employing the following *sizes for the randomly selected subsets of the data*: Y_{m0} : 1000 USPS/NDC points, (500 testing + 500 training), Y_{m1} : 2000 USPS/NDC points, (1000 testing + 1000 training), Y_{m2} : 4000 USPS/NDC points, (2000 testing + 2000 training) and Y_{m3} : 8000 USPS/NDC points, (4000 testing + 4000 training).

A point in Y_{mi} is termed a *support vector* if the dual variable of its classification constraint is nonzero at the solution, implying its classification constraint is active at the solution. Therefore, for similar testing correctness, the number of support vectors shown below is a measure of the robustness of the classifier in the sense that the support vectors include (in addition to "error" points) those points that are correctly classified, but lie in the margin and hence are not "widely separated" from the points of the other category.

The following implementation choices were made:

- $\alpha=1.25$.
- The initial value set Z_0 for the initial approximation problem is a random subset of $t Y_{m0}$ points with $u \neq 0$ after step (a) in stage 0. The w_0 are their u multipliers.

- For $i \geq 0$, the initial value set Z_{i+1} is the best Z_i^* obtained from an approximation problem in step (c), over p runs, in terms of testing set classification in step (c). The initial w_{i+1} are the corresponding w , also after step (c).

In order to evaluate the efficiency of the CCS algorithm, it is instructive to compare the quality of the solution at step (d) of stage i with the output of the **standard 1-norm classification problem** in which $C_i = Y_{mi}$:

$$\min_{u,\nu,E} \quad \nu \|E\|_1 + \|u\|_1 \text{ s.t. } D(uK(Y_{mi}, Y_{mi}) - \gamma \mathbf{1}) + E \geq \mathbf{1}, E \geq \mathbf{0}. \quad (3)$$

Note that, in the particular case of stage 0, the problem (3) is the same as step (b). However, relative to stages 1, 2 and 3, the number of classifier points that are allowed and the corresponding number of variables grow significantly in the case of problem (3).

The results for **stage 0** were:

- # of classifier points, p runs: 96.8 (step (b), avg.), 10 (step (d), max.).
- Avg. training correctness: 84.3% (step (b)), 84.38% (step (d)).
- Avg. testing correctness: 72.12% (step (b)), 72.42% (step (d)).
- Avg. # of support vectors out of 500: 362.5 (step (b)), 212.5 (step (d)).

We see that the number of support vectors in step (d) is considerably less than in the standard 1-norm classification problem. Also, using only a small number of synthetic points does not degrade testing correctness, but instead yields a small improvement (similar results are obtained for the larger subsets considered below).

Step (c) yields an approximation of the classifier obtained in step (b). In order to evaluate the need of step (d), we evaluated the average quality of the step (c) classifier over all runs in stage 0 (“pure” synthetic-point classifier, without the re-classification step (d)): we observed an average testing correctness of 56.14% and an average training correctness of 56.28% in step (c).

Consequently, the benefit of the re-classification in step (d) is evident. Another question that may arise is the quality of trivial alternatives to the NLP synthetic point generation procedure of step (c). Thus, we considered a simple choice of classifier points: given the subset of C_0 such that $u \neq 0$ after step (b), we took a t -point subset whose u multipliers had largest absolute value, and used them for re-classification of the training points in stage 0, yielding an average training correctness of 60.54% in modified step (d) and an average testing correctness of 58.16% using the approximate classifier from step (d).

We observe again the considerable advantage of using synthetic points. Similar results were obtained when t -point random samples of Y_{m0} were used.

In stages 1 and 2 the number of classifier points in step (b) is progressively lower. Specifically, in stage 2 we obtained the following results:

- # of classifier points, p runs: 13.1 (step (b), avg.), 10 (step (d), max.).
- Avg. training correctness: 76.77% (step (b)), 76.71% (step (d)).
- Avg. testing correctness: 76.76% (step (b)), 76.62% (step (d)).
- Avg. # of support vectors out of 2000: 1104 (step (d)).

Nevertheless, when we used the **standard 1-norm formulation** in problem (3) the results were much worse:

- Avg. # of classifier points, p runs: 122.2 (Y_{m1}), 180.6 (Y_{m2}).
- Avg. training correctness: 82.06% (Y_{m1}), 81.94% (Y_{m2}).
- Avg. testing correctness: 73.3% (Y_{m1}), 75.47% (Y_{m2}).
- Avg. # of support vectors: 672.5 out of 1000 (Y_{m1}), 1181 out of 2000 (Y_{m2}).

Note that while *training correctness* is improved by allowing all 2000 points in the stage 2 training subset to be used to construct the classifier, the resulting *testing correctness* of 75.47% is actually worse than that obtained by using the classifier points from the much smaller synthetic sets C_i (whose 100 points yield 76.76% classification) or Z_i^* (whose 10 points yield 76.62%). This apparently surprising result that we already observed in stage 0 is analogous to results in [4], where random 1%-5% subsets of the Adult Dataset [5] were allowed in the Gaussian kernel classifier. Observe that if the solution quality is measured by the number of support vectors, then by this measure as well, the smaller classifier sets provide better quality solutions than the full training set. Similar results are obtained when the full dataset is used in stage 3, except that the standard 1-norm LP could not be solved in 48 hours, illustrating the scalability difficulties associated with the standard approach.

We thus compared our results (as provided by the use of synthetic points in the classifier) with *random subset classifiers*. Let θ be the number of random points chosen in Y_{m2} . We observed the following:

- Avg. # of classifier points, 1-norm, p runs: 63.9 ($\theta=100$), 10 ($\theta=10$).
- Avg. training correctness: 74.28% ($\theta=100$), 59.18% ($\theta=10$).
- Avg. testing correctness: 70.83% ($\theta=100$), 57.74% ($\theta=10$).
- Avg. # of support vectors: 1814.3 out of 2000 ($\theta=10$).

For this dataset, randomly chosen subsets of the training set produce average classifications that are significantly worse than those obtained with either the 1-norm classification with the full training set or our small sets of synthetic points.

4 Conclusions and Directions for Future Research

In this paper, we have analyzed the fusion of synthetic point generators (small nonlinear programs) with chunking algorithms to develop a chunking-coordinated-synthetic (CCS) approach that achieves good generalizability while greatly reducing the size of the optimization problems that are required to produce nonlinear classifiers. Our numerical results on a modified USPS dataset show that the classifiers obtained using very small numbers of synthetic points (*as few as 10*) not only yield good generalizability (in terms of good testing set classification in ten-fold cross-validation), but also, for the noisy data considered, actually yield classifiers with *better generalizability* than either various other choices of reduced sets of training points or even the full training set (the

latter appears to result in over-training as was noted in [4] for other datasets and other reduced set approaches). Finally, since the CCS approach utilizes the solution of independent optimization problems at each stage, computation may be further accelerated by parallelization. One of our future directions for research will be the parallel implementation of CCS.

Finally, although the emphasis of this research is on nonlinear kernels, similar ideas could be applied to linear kernels. For them, problem (2) is trivial, since f is the identity function and (2) is solved by taking $w = 1$ and $Z = u^*C$. The interesting aspect of CCS in this case is the coordination step in which these single-point “ideal” classifiers for subsets are combined to produce a small classifier set for a larger subset. The key issue is whether these small classifier sets will produce good classifiers.

References

1. P.S. Bradley and O.L. Mangasarian. “Feature selection via concave minimization and support vector machines”. In J. Shavlik, editor, *Machine Learning Proceedings of the Fifteenth International Conference (ICML’98)*, pp. 82-90, San Francisco CA, 1998, Morgan Kaufmann.
2. C. J. C. Burges. “Simplified Support Vector Decision Rules”. In L. Saitta, editor, *Proceedings 13th Intl. Conf. on Machine Learning*, pp. 71-77, San Mateo CA, 1996, Morgan Kaufmann.
3. C. J. C. Burges and B. Scholkopf. “Improving the Accuracy and Speed of Support Vector Machines”. In M. Mozer, M. Jordan, and T. Petsche, editors, *Advances in Neural Information Processing Systems 9*, pages 375-381, Cambridge, MA, 1997. MIT Press.
4. Y.-J. Lee and O. L. Mangasarian, “RSVM: Reduced Support Vector Machines”. Data Mining Institute technical report 00-07, July 2000.
5. P. M. Murphy and A. W. Aha. UCI repository of machine learning databases, 1992. www.ics.uci.edu/~mllearn/MLRepository.html.
6. D. R. Musicant. NDC: Normally Distributed Clustered datasets, 1998. www.cs.wisc.edu/~musicant/data/ndc.
7. E. Osuna and F. Girosi. “Reducing the Run-time complexity of Support Vector Machines”. In *Proceedings of the 14th International Conference on Pattern Recognition*, Brisbane, Australia, 1998.
8. B. Scholkopf, S. Mika, C. J. C. Burges, P. Knirsch, K.-R. Muller, G. Ratsch and A. J. Smola. “Input Space vs. Feature Space in Kernel-Based Methods”. *IEEE Transactions on Neural Networks* 10(5):1000-1017, 1999.
9. B. Scholkopf and A. J. Smola. Kernel machines page, 2000. www.kernel-machines.org.
10. D. DeCoste and B. Scholkopf. “Training Invariant Support Vector Machines”. *Machine Learning* 46: 161-190, 2002.